

University of Groningen

MOLGENIS Research

van der Velde, K Joeri; Imhann, Floris; Charbon, Bart; Pang, Chao; van Enckevort, David; Slofstra, Mariska; Barbieri, Ruggero; Alberts, Rudi; Hendriksen, Dennis; Kelpin, Fleur

Published in:
Bioinformatics (Oxford, England)

DOI:
[10.1093/bioinformatics/bty742](https://doi.org/10.1093/bioinformatics/bty742)

IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.

Document Version
Publisher's PDF, also known as Version of record

Publication date:
2019

[Link to publication in University of Groningen/UMCG research database](#)

Citation for published version (APA):

van der Velde, K. J., Imhann, F., Charbon, B., Pang, C., van Enckevort, D., Slofstra, M., Barbieri, R., Alberts, R., Hendriksen, D., Kelpin, F., de Haan, M., de Boer, T., Haakma, S., Stroomberg, C., Scholtens, S., van de Geijn, G.-J., Festen, E. A. M., Weersma, R. K., & Swertz, M. A. (2019). MOLGENIS Research: Advanced bioinformatics data software for non-bioinformaticians. *Bioinformatics (Oxford, England)*, 35(6), 1076-1078. <https://doi.org/10.1093/bioinformatics/bty742>

Copyright

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).

The publication may also be distributed here under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license. More information can be found on the University of Groningen website: <https://www.rug.nl/library/open-access/self-archiving-pure/taverne-amendment>.

Take-down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.

Applications note

MOLGENIS Research: Advanced bioinformatics data software for non-bioinformaticians

K. Joeri van der Velde^{1,2}, Floris Imhann^{2,3}, Bart Charbon¹, Chao Pang¹, David van Enckevort¹, Mariska Slofstra¹, Ruggero Barbieri^{2,3}, Rudi Alberts^{2,3}, Dennis Hendriksen¹, Fleur Kelpin¹, Mark de Haan¹, Tommy de Boer¹, Sido Haakma¹, Connor Stroomberg¹, Salome Scholtens¹, Gert-Jan van de Geijn¹, Eleonora A. M. Festen^{2,3}, Rinse K. Weersma³ and Morris A. Swertz^{1,2,*}

¹University of Groningen and University Medical Center Groningen, Genomics Coordination Center, Groningen, The Netherlands, ²University of Groningen and University Medical Center Groningen, Department of Genetics, Groningen, The Netherlands, ³University of Groningen and University Medical Center Groningen, Department of Gastroenterology and Hepatology, Groningen, The Netherlands.

*To whom correspondence should be addressed.

Associate Editor: XXXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Motivation: The volume and complexity of biological data increases rapidly. Many clinical professionals and biomedical researchers without a bioinformatics background are generating big 'omics' data, but do not always have the tools to manage, process or publicly share these data.

Results: Here we present MOLGENIS Research, an open-source web-application to collect, manage, analyze, visualize and share large and complex biomedical data sets, without the need for advanced bioinformatics skills.

Availability and implementation: MOLGENIS Research is freely available (open source software). It can be installed from source code (see <http://github.com/molgenis>), downloaded as a precompiled WAR file (for your own server), setup inside a Docker container (see <http://molgenis.github.io>), or requested as a Software-as-a-Service subscription. For a public demo instance and complete installation instructions see <http://molgenis.org/research>.

Contact: m.a.swertz@rug.nl

1 Introduction

In order to improve human health, biomedical scientists are increasingly using large and complex data sets to discover biological mechanisms. Large numbers of patients and control participants are screened with questionnaires, biomedical measurements, high-throughput techniques such as next-generation sequencing of the genome, the transcriptome and the microbiome (Ginsburg, 2014), resulting in large quantities of phenotypic and molecular data (Bowdin *et al.*, 2014). However, many clinical professionals and biomedical researchers do not always have the proper tools to process, manage, analyze, visualize and publicly share these data (Jagadish, 2015) while complying to 'FAIR' (Findable, Accessible, In-

teroperable, and Reusable) (Wilkinson *et al.*, 2016) and 'ELSI' (Ethical, Legal and Social Implications) principles.

Several challenges arise when developing software for big data used by biomedical researchers (Raghupathi and Raghupathi, 2014). The first challenge is data capture and data management. Data systems need to be adaptable enough to not only handle today's data, but also be able seamlessly capture tomorrow's data formats (M A Swertz *et al.*, 2010). Current systems are often too strict in terms of importing new data types. As a consequence, systems must sometimes even be taken offline for database redesign (Morris A Swertz *et al.*, 2010; Adamusiak *et al.*, 2012). Therefore, a good system needs to allow continuous use while databases can be redesigned and unforeseen data types can be

and persistence, indexing its data tables, offering HTTP access and authorization, and tools to connect data to FAIR vocabularies such as ontologies (Pang *et al.*, 2014). For joint analysis of data sets, we have developed *Mapping Service* tools to make both columns (Pang *et al.*, 2014) and values (Pang *et al.*, 2015) interoperable between data sets so they can be merged. The *Tag Wizard* app can assign meaning to data columns using ontologies, which can be integrated across different data sets using the *Mapping Service* app. Data sets and variables can be made findable without exposing (sensitive) data values by creating a catalogue from a combination of raw data, curated data or interesting results collected in the system. Others can browse this catalogue before contacting or submitting a request for access. MOLGENIS Research supports the complete data access and request workflow designed by the data owner. Super users can also create FAIR endpoints (Wilkinson, Verborgh, *et al.*, 2017) based on definitions of Metadata, Catalog, Dataset, Distribution and Response, which ensures your data is machine-findable and thereby has increased findability.

3 Implementation

MOLGENIS Research is implemented using open and freely usable industry standards. It is available under the GNU Lesser General Public License v3.0 (<https://www.gnu.org/licenses/lgpl-3.0.en.html>). It is written in Java 1.8 (<https://java.com>), supported by the Spring MVC framework (<https://spring.io>). It uses Apache Maven (<https://maven.apache.org>) to manage dependencies, and runs on an Apache Tomcat (<http://tomcat.apache.org>) webserver. Data is stored in a PostgreSQL database (<https://www.postgresql.org>) and indexed by Elasticsearch (<https://www.elastic.co>) for high performance and horizontal scaling ability by data replication and sharding, respectively. Final storage and query performance depends on specific hardware and software configuration. Its graphical user interface is composed of Bootstrap (<https://getbootstrap.com>), Vue (<https://vuejs.org>) and Freemarker templates (<https://freemarker.apache.org>). FAIR endpoints are implemented in W3C RDF 1.1 Turtle (<https://www.w3.org/TR/turtle>).

4 Conclusion

We have built MOLGENIS Research, a web application for the biomedical field to work with multi-omics data sets without being dependent on bioinformaticians. MOLGENIS Research enables researchers to more efficiently collect, manage, analyze, visualize and share data, as well as offering support to make data FAIR in a flexible and safe way. MOLGENIS Research offers all the advantages of a true database system with detailed data management and access control options, while at the same time being able to grow ‘organically’ by allowing data to be dynamically shaped based on what is needed in practice, and adding custom extensions such as visualizations and algorithms into a running system without downtime. It can be used as a project database from day one as there is no need to design a data model upfront.

Currently, MOLGENIS Research has been adopted by several research projects, including 1000IBD, 500FG and LifeLines. The 1000IBD database (<http://1000ibd.org>) contains a range of clinical and research phenotypes for up to 2,000 patients per -omics type, which includes quantifications of 12,000+ microbiome OTUs, 400+ immuno-chip markers, and ~300 RNA-seq experiments. The 500FG database (<https://hfgp.bbmri.nl>) contains microbiome, metabolomics, cytokine, QTL, cell staining, serum Ig and flow cytometry data for around 500 individuals. Identifier codes for individuals serve as foreign keys that can

link data tables together for data integration and analysis. Lastly, the LifeLines data catalogue (<https://catalogue.lifelines.nl>) contains the metadata for around 40,000 data items available for researchers such as questionnaires, measurements and (blood and urine) sample analyses from a longitudinal study of 167,000 individuals. We expect more projects to follow soon, and gladly invite everyone to help us in expanding and evolving the MOLGENIS Research solution to serve all popular research needs. We strongly encourage interested users to try the demo, download and install MOLGENIS Research at <http://molgenis.org/research>.

Acknowledgements

We thank Benjamin Kant for feedback and comments.

Funding

We thank BBMRI-NL for sponsoring the development of the software described in this manuscript via a voucher. BBMRI-NL is a research infrastructure financed by the Netherlands Organization for Scientific Research (NWO), grant number 184.033.111. We also thank NWO VIDI grant number 917.164.455.

Conflict of Interest: none declared.

References

- Adamusiak, T. *et al.* (2012) Observ-OM and Observ-TAB: Universal syntax solutions for the integration, search and exchange of phenotype and genotype information. *Hum. Mutat.*, **33**, 867–73.
- Auffray, C. *et al.* (2016) Making sense of big data in health research: Towards an EU action plan. *Genome Med*, **8**.
- Bowdin, S. *et al.* (2014) The Genome Clinic: A Multidisciplinary Approach to Assessing the Opportunities and Challenges of Integrating Genomic Analysis into Clinical Care. *Hum. Mutat.*, **35**, 513–519.
- Down, T.A. *et al.* (2011) Dalliace: interactive genome viewing on the web. *Bioinformatics*, **27**, 889–890.
- Ginsburg, G. (2014) Medical genomics: Gather and use genetic data in health care. *Nature*, **508**, 451–453.
- Higdon, R. *et al.* (2015) The Promise of Multi-Omics and Clinical Data Integration to Identify and Target Personalized Healthcare Approaches in Autism Spectrum Disorders. *Omi. A J. Integr. Biol.*, **19**, 197–208.
- Jagadish, H. V (2015) Big Data and Science: Myths and Reality. *Big Data Res.*, **2**, 49–52.
- Pang, C. *et al.* (2014) BiobankConnect: software to rapidly connect data elements for pooled analysis across biobanks using ontological and lexical indexing. *J. Am. Med. Informatics Assoc.*, 65–75.
- Pang, C. *et al.* (2015) SORTA: A system for ontology-based re-coding and technical annotation of biomedical phenotype data. *Database*, **2015**.
- Raghupathi, W. and Raghupathi, V. (2014) Big data analytics in healthcare: promise and potential. *Heal. Inf. Sci. Syst.*, **2**, 3.
- Stieb, D.M. *et al.* (2017) Promise and pitfalls in the application of big data to occupational and environmental health. *BMC Public Health*, **17**.
- Suravajhala, P. *et al.* (2016) Multi-omic data integration and analysis using systems genomics approaches: methods and applications in animal production, health and welfare. *Genet. Sel. Evol.*, **48**.
- Swertz, M.A. *et al.* (2010) The MOLGENIS toolkit: rapid prototyping of

-
- bioinformatics at the push of a button. *BMC Bioinformatics*, **11**.
- Swertz,M.A. *et al.* (2010) XGAP: a uniform and extensible data model and software platform for genotype and phenotype experiments. *Genome Biol.*, **11**, R27.
- van der Velde,K.J. *et al.* (2017) GAVIN: Gene-Aware Variant Interpretation for medical sequencing. *Genome Biol.*, **18**, 6.
- Wilkinson,M.D., Sansone,S.-A., *et al.* (2017) A design framework and exemplar metrics for FAIRness. *bioRxiv*.
- Wilkinson,M.D., Verborgh,R., *et al.* (2017) Interoperability and FAIRness through a novel combination of Web technologies. *PeerJ Comput. Sci.*, **3**, e110.
- Wilkinson,M.D. *et al.* (2016) The FAIR Guiding Principles for scientific data management and stewardship. *Sci. Data*, **3**, 160018.